
CHAPTER 8

Hidden Markov Models

The hidden Markov model (HMM) is a very powerful statistical method of characterizing the observed data samples of a discrete-time series. Not only can it provide an efficient way to build parsimonious parametric models, but can also incorporate the dynamic programming principle in its core for a unified pattern segmentation and pattern classification of time-varying data sequences. The data samples in the time series can be discretely or continuously distributed; they can be scalars or vectors. The underlying assumption of the HMM is that the data samples can be well characterized as a parametric random process, and the parameters of the stochastic process can be estimated in a precise and well-defined framework. The basic HMM theory was published in a series of classic papers by Baum and his colleagues [4]. The HMM has become one of the most powerful statistical methods for modeling speech signals. Its principles have been successfully used in automatic speech recognition, formant and pitch tracking, speech enhancement, speech synthesis, statistical language modeling, part-of-speech tagging, spoken language understanding, and machine translation [3, 4, 8, 10, 12, 18, 23, 37].

8.1. THE MARKOV CHAIN

A *Markov chain* models a class of random processes that incorporates a minimum amount of memory without being completely memoryless. In this subsection we focus on the discrete-time Markov chain only.

Let $\mathbf{X} = X_1, X_2, \dots, X_n$ be a sequence of random variables from a finite discrete alphabet $O = \{o_1, o_2, \dots, o_M\}$. Based on the Bayes rule, we have

$$P(X_1, X_2, \dots, X_n) = P(X_1) \prod_{i=2}^n P(X_i | X_1^{i-1}) \quad (8.1)$$

where $X_1^{i-1} = X_1, X_2, \dots, X_{i-1}$. The random variables X are said to form a first-order Markov chain, provided that

$$P(X_i | X_1^{i-1}) = P(X_i | X_{i-1}) \quad (8.2)$$

As a consequence, for the first-order Markov chain, Eq. (8.1) becomes

$$P(X_1, X_2, \dots, X_n) = P(X_1) \prod_{i=2}^n P(X_i | X_{i-1}) \quad (8.3)$$

Equation (8.2) is also known as the *Markov assumption*. This assumption uses very little memory to model dynamic data sequences: the probability of the random variable at a given time depends only on the value at the preceding time. The Markov chain can be used to model time-invariant (stationary) events if we discard the time index i ,

$$P(X_i = s | X_{i-1} = s') = P(s | s') \quad (8.4)$$

If we associate X_i to a state, the Markov chain can be represented by a finite state process with transition between states specified by the probability function $P(s | s')$. Using this finite state representation, the Markov assumption (Eq. (8.2)) is translated to the following: the probability that the Markov chain will be in a particular state at a given time depends only on the state of the Markov chain at the previous time.

Consider a Markov chain with N distinct states labeled by $\{1, \dots, N\}$, with the state at time t in the Markov chain denoted as s_t ; the parameters of a Markov chain can be described as follows:

$$a_{ij} = P(s_t = j | s_{t-1} = i) \quad 1 \leq i, j \leq N \quad (8.5)$$

$$\pi_i = P(s_1 = i) \quad 1 \leq i \leq N \quad (8.6)$$

where a_{ij} is the transition probability from state i to state j ; and π_i is the initial probability that the Markov chain will start in state i . Both transition and initial probabilities are bound to the constraints:

$$\sum_{j=1}^N a_{ij} = 1; \quad 1 \leq i \leq N$$

$$\sum_{j=1}^N \pi_j = 1$$
(8.7)

The Markov chain described above is also called the observable Markov model because the output of the process is the set of states at each time instance t , where each state corresponds to an observable event X_t . In other words, there is one-to-one correspondence between the observable event sequence \mathbf{X} and the Markov chain state sequence $\mathbf{S} = s_1, s_2, \dots, s_n$. Consider a simple three-state Markov chain for the Dow Jones Industrial average as shown in Figure 8.1. At the end of each day, the Dow Jones Industrial average may correspond to one of the following states:

- state 1** – *up* (in comparison to the index of previous day)
- state 2** – *down* (in comparison to the index of previous day)
- state 3** – *unchanged* (in comparison to the index of previous day)

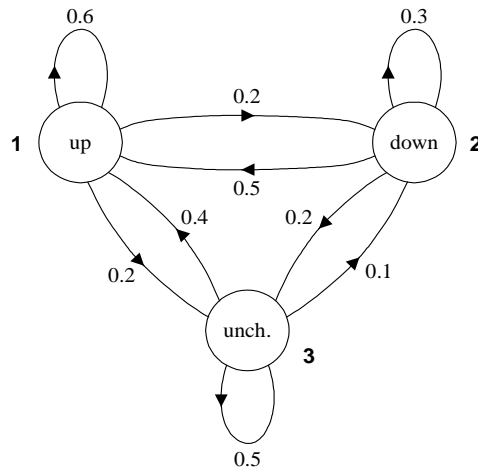


Figure 8.1 A Markov chain for the Dow Jones Industrial average. Three states represent *up*, *down*, and *unchanged* respectively.

The parameter for this Dow Jones Markov chain may include a state-transition probability matrix

$$A = \{a_{ij}\} = \begin{bmatrix} 0.6 & 0.2 & 0.2 \\ 0.5 & 0.3 & 0.2 \\ 0.4 & 0.1 & 0.5 \end{bmatrix}$$

and an initial state probability matrix

$$\boldsymbol{\pi} = (\pi_i)^t = \begin{pmatrix} 0.5 \\ 0.2 \\ 0.3 \end{pmatrix}$$

Suppose you would like to find out the probability for five consecutive *up* days. Since the observed sequence of *up-up-up-up-up* corresponds to the state sequence of (1, 1, 1, 1, 1), the probability will be

$$\begin{aligned} P(5 \text{ consecutive } up \text{ days}) &= P(1, 1, 1, 1, 1) \\ &= \pi_1 a_{11} a_{11} a_{11} a_{11} = 0.5 \times (0.6)^4 = 0.0648 \end{aligned}$$

8.2. DEFINITION OF THE HIDDEN MARKOV MODEL

In the Markov chain, each state corresponds to a deterministically observable event; i.e., the output of such sources in any given state is not random. A natural extension to the Markov chain introduces a non-deterministic process that generates output observation symbols in any given state. Thus, the observation is a probabilistic function of the state. This new model is known as a *hidden Markov model*, which can be viewed as a double-embedded stochastic process with an underlying stochastic process (the state sequence) not directly observable. This underlying process can only be probabilistically associated with another observable stochastic process producing the sequence of features we can observe.

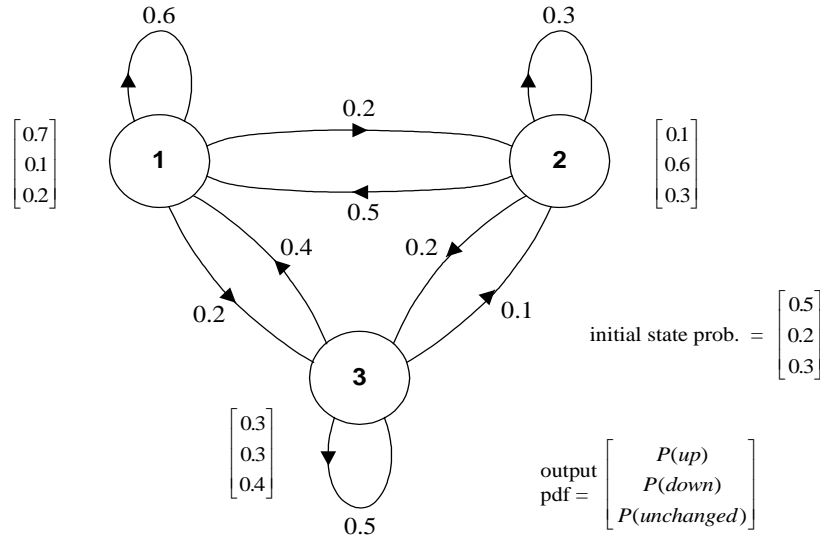


Figure 8.2 A hidden Markov model for the Dow Jones Industrial average. The three states no longer have deterministic meanings as the Markov chain illustrated in Figure 8.1.

A hidden Markov model is basically a Markov chain where the output observation is a random variable X generated according to an output probabilistic function associated with

each state. Figure 8.2 shows a revised hidden Markov model for the Dow Jones Industrial average. You see that each state now can generate all three output observations: *up*, *down*, and *unchanged* according to its output pdf. This means that there is no longer a one-to-one correspondence between the observation sequence and the state sequence, so you cannot unanimously determine the state sequence for a given observation sequence; i.e., the state sequence is not observable and therefore hidden. This is why the world *hidden* is placed in front of *Markov models*. Although the state of an HMM is hidden, it often contains salient information about the data we are modeling. For example, the first state in Figure 8.2 indicates a *bull* market, and the second state indicates a *bear* market as specified by the output probability in each state. Formally speaking, a hidden Markov model is defined by:

- $O = \{o_1, o_2, \dots, o_M\}$ - An output observation alphabet.¹ The observation symbols correspond to the physical output of the system being modeled. In the case of the Dow Jones Industrial average HMM, the output observation alphabet is the collection of three categories - $O = \{up, down, unchanged\}$.
- $\Omega = \{1, 2, \dots, N\}$ - A set of states representing the state space. Here s_t is denoted as the state at time t . In the case of the Dow Jones Industrial average HMM, the state may indicate a *bull* market, a *bear* market, and a *stable* market.
- $A = \{a_{ij}\}$ - A transition probability matrix, where a_{ij} is the probability of taking a transition from state i to state j , i.e.,

$$a_{ij} = P(s_t = j | s_{t-1} = i) \quad (8.8)$$

- $B = \{b_i(k)\}$ - An output probability matrix,² where $b_i(k)$ is the probability of emitting symbol o_k when state i is entered. Let $\mathbf{X} = X_1, X_2, \dots, X_t, \dots$ be the observed output of the HMM. The state sequence $S = s_1, s_2, \dots, s_t, \dots$ is not observed (hidden), and $b_i(k)$ can be rewritten as follows:

$$b_i(k) = P(X_t = o_k | s_t = i) \quad (8.9)$$

- $\pi = \{\pi_i\}$ - A initial state distribution where

$$\pi_i = P(s_0 = i) \quad 1 \leq i \leq N \quad (8.10)$$

Since a_{ij} , $b_{ij}(k)$, and π_i are all probabilities, they must satisfy the following properties:

¹ Although we use the discrete output observation here, we can extend it to the continuous case with a continuous pdf. You can also use vector quantization to map a continuous vector variable into a discrete alphabet set.

² The output distribution can also be transition-dependent. Although these two formations look different, the state-dependent one can be reformulated as a transition-dependent one with the constraint of all the transitions entering the same state sharing the same output distribution.

$$a_{ij} \geq 0, \quad b_i(k) \geq 0, \quad \pi_i \geq 0 \quad \forall \text{ all } i, j, k \quad (8.11)$$

$$\sum_{j=1}^N a_{ij} = 1 \quad (8.12)$$

$$\sum_{k=1}^M b_i(k) = 1 \quad (8.13)$$

$$\sum_{i=1}^N \pi_i = 1 \quad (8.14)$$

To sum up, a complete specification of an HMM includes two constant-size parameters, N and M , representing the total number of states and the size of observation alphabets, observation alphabet O , and three sets (matrices) of probability measures \mathbf{A} , \mathbf{B} , $\boldsymbol{\pi}$. For convenience, we use the following notation

$$\Phi = (\mathbf{A}, \mathbf{B}, \boldsymbol{\pi}) \quad (8.15)$$

to indicate the whole parameter set of an HMM and sometimes use the parameter set Φ to represent the HMM interchangeably without ambiguity.

In the first-order hidden Markov model discussed above, there are two assumptions. The first is the *Markov assumption* for the Markov chain.

$$P(s_t | s_1^{t-1}) = P(s_t | s_{t-1}) \quad (8.16)$$

where s_1^{t-1} represents the state sequence s_1, s_2, \dots, s_{t-1} . The second is the *output-independence assumption*:

$$P(X_t | X_1^{t-1}, s_1^t) = P(X_t | s_t) \quad (8.17)$$

where X_1^{t-1} represents the output sequence X_1, X_2, \dots, X_{t-1} . The output-independence assumption states that the probability that a particular symbol is emitted at time t depends only on the state s_t and is conditionally independent of the past observations.

You might argue that these assumptions limit the memory of the first-order hidden Markov models and may lead to model deficiency. However, in practice, they make evaluation, decoding, and learning feasible and efficient without significantly affecting the modeling capability, since those assumptions greatly reduce the number of parameters that need to be estimated.

Given the definition of HMMs above, three basic problems of interest must be addressed before they can be applied to real-world applications.

1. **The Evaluation Problem** – Given a model Φ and a sequence of observations $\mathbf{X} = (X_1, X_2, \dots, X_T)$, what is the probability $P(\mathbf{X}|\Phi)$; i.e., the probability of the model that generates the observations?
2. **The Decoding Problem** – Given a model Φ and a sequence of observations $\mathbf{X} = (X_1, X_2, \dots, X_T)$, what is the most likely state sequence $\mathbf{S} = (s_0, s_1, s_2, \dots, s_T)$ in the model that produces the observations?
3. **The Learning Problem** – Given a model Φ and a set of observations, how can we adjust the model parameter $\hat{\Phi}$ to maximize the joint probability (likelihood) $\prod_{\mathbf{X}} P(\mathbf{X}|\Phi)$?

If we could solve the *evaluation* problem, we would have a way of evaluating how well a given HMM matches a given observation sequence. Therefore, we could use HMM to do pattern recognition, since the likelihood $P(\mathbf{X}|\Phi)$ can be used to compute posterior probability $P(\Phi|\mathbf{X})$, and the HMM with highest posterior probability can be determined as the desired pattern for the observation sequence. If we could solve the decoding problem, we could find the best matching state sequence given an observation sequence, or, in other words, we could uncover the hidden state sequence. As discussed in Chapters 12 and 13, these are the basics for the decoding in continuous speech recognition. Last but not least, if we could solve the learning problem, we would have the means to automatically estimate the model parameter Φ from an ensemble of training data. These three problems are tightly linked under the same probabilistic framework. The efficient implementation of these algorithms shares the principle of dynamic programming that we briefly discuss next.

8.2.1. Dynamic Programming and DTW

The dynamic programming concept, also known as *dynamic time warping* (DTW) in speech recognition [40], has been widely used to derive the overall distortion between two speech templates. In these template-based systems, each speech template consists of a sequence of speech vectors. The overall distortion measure is computed from the accumulated distance between two feature vectors that are aligned between two speech templates with minimal overall distortion. The DTW method can warp two speech templates $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$ and $(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_M)$ in the time dimension to alleviate nonlinear distortion as illustrated in Figure 8.3.

This is roughly equivalent to the problem of finding the minimum distance in the trellis between these two templates. Associated with every pair (i, j) is a distance $d(i, j)$ between two speech vectors \mathbf{x}_i and \mathbf{y}_j . To find the optimal path between starting point $(1, 1)$ and end point (N, M) from left to right, we need to compute the optimal accumulated distance $D(N, M)$. We can enumerate all possible accumulated distance from $(1, 1)$ to (N, M) and identify the one that has the minimum distance. Since there are M possible moves for

each step from left to right in Figure 8.3, all the possible paths from $(1, 1)$ to (N, M) will be exponential. Dynamic programming principle can drastically reduce the amount of computation by avoiding the enumeration of sequences that cannot possibly be optimal. Since the same optimal path after each step must be based on the previous step, the minimum distance $D(i, j)$ must satisfy the following equation:

$$D(i, j) = \min_k [D(i-1, k) + d(k, j)] \quad (8.18)$$

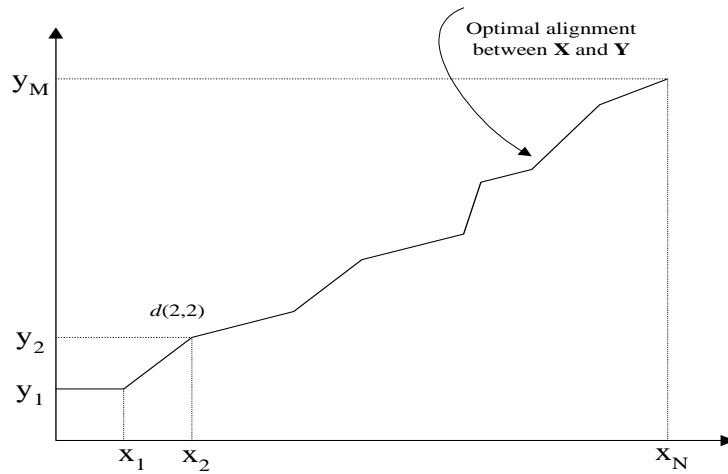


Figure 8.3 Direct comparison between two speech templates $\mathbf{X} = x_1 x_2 \dots x_N$ and $\mathbf{Y} = y_1 y_2 \dots y_M$.

Equation (8.18) indicates you only need to consider and keep only the best move for each pair although there are M possible moves. The recursion allows the optimal path search to be conducted incrementally from left to right. In essence, dynamic programming delegates the solution recursively to its own sub-problem. The computation proceeds from the small sub-problem $(D(i-1, k))$ to the larger sub-problem $(D(i, j))$. We can identify the optimal match y_j with respect to x_i and save the index in a back pointer table $B(i, j)$ as we move forward. The optimal path can be back traced after the optimal path is identified. The algorithm is described in Algorithm 8.1.

The advantage of the dynamic programming lies in the fact that once a sub-problem is solved, the partial result can be stored and never needs to be recalculated. This is a very important principle that you see again and again in building practical spoken language systems.

Speech recognition based on DTW is simple to implement and fairly effective for small-vocabulary speech recognition. Dynamic programming can temporally align patterns to account for differences in speaking rates across talkers as well as across repetitions of the word by the same talker. However, it does not have a principled way to derive an averaged template for each pattern from a large amount of training samples. A multiplicity of refer-

ence training tokens is typically required to characterize the variation among different utterances. As such, the HMM is a much better alternative for spoken language processing.

ALGORITHM 8.1: THE DYNAMIC PROGRAMMING ALGORITHM

Step 1: Initialization

$D(1,1) = d(1,1), B(1,1) = 1$, for $j = 2, \dots, M$ compute $D(1, j) = \infty$

Step 2: Iteration

for $i = 2, \dots, N$ {

for $j = 1, \dots, M$ compute {

$$D(i, j) = \min_{1 \leq p \leq M} [D(i-1, p) + d(p, j)]$$

$$B(i, j) = \arg \min_{1 \leq p \leq M} [D(i-1, p) + d(p, j)] \quad \}} \}$$

Step 3: Backtracking and Termination

The optimal (minimum) distance is $D(N, M)$ and the optimal path is (s_1, s_2, \dots, s_N)

where $s_N = M$ and $s_i = B(i+1, s_{i+1})$, $i = N-1, N-2, \dots, 1$

8.2.2. How to Evaluate an HMM – The Forward Algorithm

To calculate the probability (likelihood) $P(\mathbf{X}|\Phi)$ of the observation sequence $\mathbf{X} = (X_1, X_2, \dots, X_T)$, given the HMM Φ , the most intuitive way is to sum up the probabilities of all possible state sequences:

$$P(\mathbf{X}|\Phi) = \sum_{\text{all } \mathbf{S}} P(\mathbf{S}|\Phi) P(\mathbf{X}|\mathbf{S}, \Phi) \quad (8.19)$$

In other words, to compute $P(\mathbf{X}|\Phi)$, we first enumerate all possible state sequences \mathbf{S} of length T , that generate observation sequence \mathbf{X} , and then sum all the probabilities. The probability of each path \mathbf{S} is the product of the state sequence probability (first factor) and the joint output probability (the second factor) along the path.

For one particular state sequence $\mathbf{S} = (s_1, s_2, \dots, s_T)$, where s_1 is the initial state, the state-sequence probability in Eq. (8.19) can be rewritten by applying Markov assumption:

$$P(\mathbf{S}|\Phi) = P(s_1|\Phi) \prod_{t=2}^T P(s_t | s_{t-1}, \Phi) = \pi_{s_1} a_{s_1 s_2} \dots a_{s_{T-1} s_T} = a_{s_0 s_1} a_{s_1 s_2} \dots a_{s_{T-1} s_T} \quad (8.20)$$

where $a_{s_0 s_1}$ denotes π_{s_1} for simplicity. For the same state sequence \mathbf{S} , the joint output probability along the path can be rewritten by applying the output-independent assumption:

$$\begin{aligned}
P(\mathbf{X} | \mathbf{S}, \Phi) &= P(X_1^T | S_1^T, \Phi) = \prod_{t=1}^T P(X_t | s_t, \Phi) \\
&= b_{s_1}(X_1) b_{s_2}(X_2) \dots b_{s_T}(X_T)
\end{aligned} \tag{8.21}$$

Substituting Eq. (8.20) and (8.21) into (8.19), we get:

$$\begin{aligned}
P(\mathbf{X} | \Phi) &= \sum_{\text{all } \mathbf{S}} P(\mathbf{S} | \Phi) P(\mathbf{X} | \mathbf{S}, \Phi) \\
&= \sum_{\text{all } \mathbf{S}} a_{s_0 s_1} b_{s_1}(X_1) a_{s_1 s_2} b_{s_2}(X_2) \dots a_{s_{T-1} s_T} b_{s_T}(X_T)
\end{aligned} \tag{8.22}$$

Equation (8.22) can be interpreted as follows. First we enumerate all possible state sequence with length $T+1$. For any given state sequence, we start from initial state s_1 with probability π_{s_1} or $a_{s_0 s_1}$. We take a transition from s_{t-1} to s_t with probability $a_{s_{t-1} s_t}$ and generate the observation X_t with probability $b_{s_t}(X_t)$ until we reach the last transition.

However, direct evaluation of Eq. (8.22) according to the interpretation above requires enumeration of $O(N^T)$ possible state sequences, which results in exponential computational complexity. Fortunately, a more efficient algorithm can be used to calculate Eq. (8.22). The trick is to store intermediate results and use them for subsequent state-sequence calculations to save computation. This algorithm is known as the *forward algorithm*.

Based on the HMM assumptions, the calculation of $P(s_t | s_{t-1}, \Phi) P(X_t | s_t, \Phi)$ involves only s_{t-1}, s_t , and X_t , so, it is possible to compute the likelihood with $P(\mathbf{X} | \Phi)$ with recursion on t . Let's define forward probability:

$$\alpha_t(i) = P(X_1^t, s_t = i | \Phi) \tag{8.23}$$

$\alpha_t(i)$ is the probability that the HMM is in state i having generated partial observation X_1^t (namely $X_1 X_2 \dots X_t$). $\alpha_t(i)$ can be calculated inductively as illustrated in Algorithm 8.2. This inductive procedure shown in Eq. (8.24) can be illustrated in a *trellis*. Figure 8.4 illustrates the computation of forward probabilities α via a trellis framework for the Dow Jones Industrial average HMM shown in Figure 8.2. Given two consecutive *up* days for the Dow Jones Industrial average, we can find the forward probability α based on the model of Figure 8.2. An arrow in Figure 8.4 indicates a transition from its origin state to its destination state. The number inside each cell indicates the forward probability α . We start the α cells from $t = 0$, where the α cells contains exactly the initial probabilities. The other cells are computed in a *time-synchronous* fashion from left to right, where each cell for time t is completely computed before proceeding to time $t+1$. When the states in the last column have been computed, the sum of all probabilities in the final column is the probability of generating the observation sequence. For most speech problems, we need to have the HMM end in some particular exit state (a.k.a final state, s_F), and we thus have $P(\mathbf{X} | \Phi) = \alpha_T(s_F)$.

It is easy to show that the complexity for the forward algorithm is $O(N^2T)$ rather than the exponential one. This is because we can make full use of partially computed probabilities for the improved efficiency.

ALGORITHM 8.2 THE FORWARD ALGORITHM

Step 1: Initialization

$$\alpha_1(i) = \pi_i b_i(X_1) \quad 1 \leq i \leq N$$

Step 2: Induction

$$\alpha_t(j) = \left[\sum_{i=1}^N \alpha_{t-1}(i) a_{ij} \right] b_j(X_t) \quad 2 \leq t \leq T; \quad 1 \leq j \leq N \quad (8.24)$$

Step 3: Termination

$$P(\mathbf{X}|\Phi) = \sum_{i=1}^N \alpha_T(i) \quad \text{If it is required to end in the final state, } P(\mathbf{X}|\Phi) = \alpha_T(s_F)$$

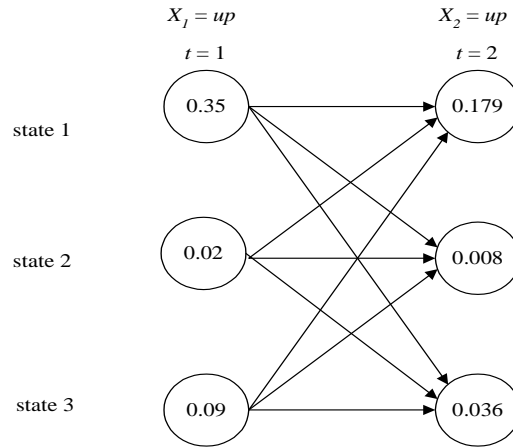


Figure 8.4 The forward trellis computation for the HMM of the Dow Jones Industrial average.

8.2.3. How to Decode an HMM - The Viterbi Algorithm

The forward algorithm, in the previous section, computes the probability that an HMM generates an observation sequence by summing up the probabilities of all possible paths, so it does not provide the best path (or state sequence). In many applications, it is desirable to find such a path. As a matter of fact, finding the best path (state sequence) is the cornerstone for searching in continuous speech recognition. Since the state sequence is hidden (unobserved) in the HMM framework, the most widely used criterion is to find the state sequence

that has the highest probability of being taken while generating the observation sequence. In other words, we are looking for the state sequence $\mathbf{S} = (s_1, s_2, \dots, s_T)$ that maximizes $P(\mathbf{S}, \mathbf{X} | \Phi)$. This problem is very similar to the optimal-path problem in dynamic programming. As a consequence, a formal technique based on dynamic programming, known as *Viterbi* algorithm [43], can be used to find the best state sequence for an HMM. In practice, the same method can be used to evaluate HMMs that offers an approximate solution close to the case obtained using the forward algorithm described in Section 8.2.2.

The Viterbi algorithm can be regarded as the dynamic programming algorithm applied to the HMM or as a modified forward algorithm. Instead of summing up probabilities from different paths coming to the same destination state, the Viterbi algorithm picks and remembers the best path. To define the best-path probability:

$$V_t(i) = P(X_1^t, S_1^{t-1}, s_t = i | \Phi) \quad (8.24)$$

$V_t(i)$ is the probability of the most likely state sequence at time t , which has generated the observation X_1^t (until time t) and ends in state i . A similar induction procedure for the Viterbi algorithm can be described in Algorithm 8.3.

ALGORITHM 8.3 THE VITERBI ALGORITHM

Step 1: Initialization

$$V_1(i) = \pi_i b_i(X_1) \quad 1 \leq i \leq N$$

$$B_1(i) = 0$$

Step 2: Induction

$$V_t(j) = \max_{1 \leq i \leq N} [V_{t-1}(i) a_{ij}] b_j(X_t) \quad 2 \leq t \leq T; \quad 1 \leq j \leq N \quad (8.25)$$

$$B_t(j) = \text{Arg max}_{1 \leq i \leq N} [V_{t-1}(i) a_{ij}] \quad 2 \leq t \leq T; \quad 1 \leq j \leq N \quad (8.26)$$

Step 3: Termination

$$\text{The best score} = \max_{1 \leq i \leq N} [V_T(i)]$$

$$s_T^* = \text{Arg max}_{1 \leq i \leq N} [B_T(i)]$$

Step 4: Backtracking

$$s_t^* = B_{t+1}(s_{t+1}^*) \quad t = T-1, T-2, \dots, 1$$

$\mathbf{S}^* = (s_1^*, s_2^*, \dots, s_T^*)$ is the best sequence

This Viterbi algorithm can also be illustrated in a trellis framework similar to the one for the forward algorithm shown in Figure 8.4. Instead of summing up all the paths, Figure 8.5 illustrates the computation of V_t by picking the best path in each cell. The number inside

each cell indicates the best score γ_t and the best path leading to each cell is indicated by solid lines while the rest of the paths are indicated by dashed line. Again, the computation is done in a *time-synchronous* fashion from left to right. The complexity for the Viterbi algorithm is also $O(N^2T)$.

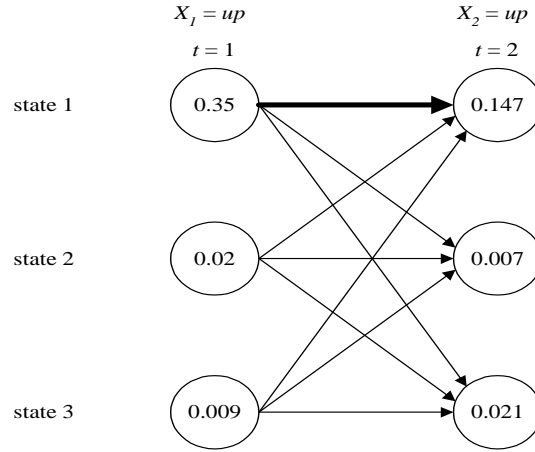


Figure 8.5 The Viterbi trellis computation for the HMM of the Dow Jones Industrial average.

8.2.4. How to Estimate HMM Parameters – Baum-Welch Algorithm

It is very important to estimate the model parameters $\Phi = (\mathbf{A}, \mathbf{B}, \boldsymbol{\pi})$ to accurately describe the observation sequences. This is by far the most difficult of the three problems, because there is no known analytical method that maximizes the joint probability of the training data in a closed form. Instead, the problem can be solved by the iterative *Baum-Welch* algorithm, also known as the *forward-backward* algorithm.

The HMM learning problem is a typical case of unsupervised learning discussed in Chapter 4, where the data is incomplete because of the hidden state sequence. The EM algorithm is perfectly suitable for this problem. As a matter of fact, Baum and colleagues used the same principle as that of the EM algorithm. Before we describe the formal Baum-Welch algorithm, we first define a few useful terms. In a manner similar to the forward probability, we define backward probability as:

$$\beta_t(i) = P(X_{t+1}^T | s_t = i, \Phi) \quad (8.27)$$

where $\beta_t(i)$ is the probability of generating partial observation X_{t+1}^T (from $t+1$ to the end) given that the HMM is in state i at time t , $\beta_t(i)$ can then be calculated inductively;

Initialization:

$$\beta_T(i) = 1/N \quad 1 \leq i \leq N$$

Induction:

$$\beta_t(i) = \left[\sum_{j=1}^N a_{ij} b_j(X_{t+1}) \beta_{t+1}(j) \right] \quad t=T-1 \dots 1; \quad 1 \leq i \leq N \quad (8.28)$$

The relationship of adjacent α and β (α_{t-1} & α_t and β_t & β_{t+1}) can be best illustrated as shown in Figure 8.6. α is computed recursively from left to right, and β recursively from right to left.

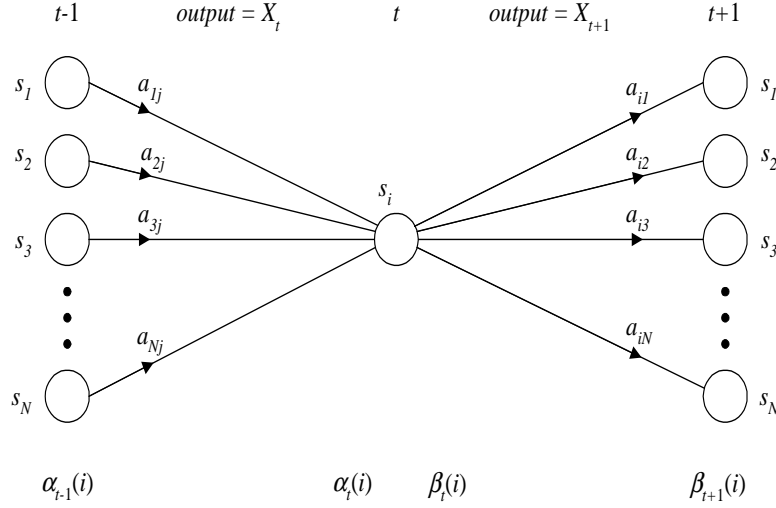


Figure 8.6 The relationship of α_{t-1} & α_t and β_t & β_{t+1} in the forward-backward algorithm.

Next, we define $\gamma_t(i, j)$, which is the probability of taking the transition from state i to state j at time t , given the model and observation sequence, i.e.

$$\begin{aligned} \gamma_t(i, j) &= P(s_{t-1} = i, s_t = j | X_1^T, \Phi) \\ &= \frac{P(s_{t-1} = i, s_t = j, X_1^T | \Phi)}{P(X_1^T | \Phi)} \\ &= \frac{\alpha_{t-1}(i) a_{ij} b_j(X_t) \beta_t(j)}{\sum_{k=1}^N \alpha_t(k)} \end{aligned} \quad (8.29)$$

The equation above can be best illustrated as shown in Figure 8.7.

We can iteratively refine the HMM parameter vector $\Phi = (\mathbf{A}, \mathbf{B}, \boldsymbol{\pi})$ by maximizing the likelihood $P(\mathbf{X} | \Phi)$ for each iteration. We use $\hat{\Phi}$ to denote the new parameter vector derived from the parameter vector Φ in the previous iteration. According to the EM algorithm

of Chapter 4, the maximization process is equivalent to maximizing the following Q -function:

$$Q(\Phi, \hat{\Phi}) = \sum_{\text{all } S} \frac{P(\mathbf{X}, \mathbf{S} | \Phi)}{P(\mathbf{X} | \Phi)} \log P(\mathbf{X}, \mathbf{S} | \hat{\Phi}) \quad (8.30)$$

where $P(\mathbf{X}, \mathbf{S} | \Phi)$ and $\log P(\mathbf{X}, \mathbf{S} | \hat{\Phi})$ can be expressed as:

$$P(\mathbf{X}, \mathbf{S} | \Phi) = \prod_{t=1}^T a_{s_{t-1}s_t} b_{s_t}(X_t) \quad (8.31)$$

$$\log P(\mathbf{X}, \mathbf{S} | \Phi) = \sum_{t=1}^T \log a_{s_{t-1}s_t} + \sum_{t=1}^T \log b_{s_t}(X_t) \quad (8.32)$$

Equation (8.30) can thus be rewritten as

$$Q(\Phi, \hat{\Phi}) = Q_{a_i}(\Phi, \hat{a}_i) + Q_{b_j}(\Phi, \hat{b}_j) \quad (8.33)$$

where

$$Q_{a_i}(\Phi, \hat{a}_i) = \sum_i \sum_j \sum_t \frac{P(\mathbf{X}, s_{t-1} = i, s_t = j | \Phi)}{P(\mathbf{X} | \Phi)} \log \hat{a}_{ij} \quad (8.34)$$

$$Q_{b_j}(\Phi, \hat{b}_j) = \sum_j \sum_k \sum_{t \in X_t = o_k} \frac{P(\mathbf{X}, s_t = j | \Phi)}{P(\mathbf{X} | \Phi)} \log \hat{b}_j(k) \quad (8.35)$$

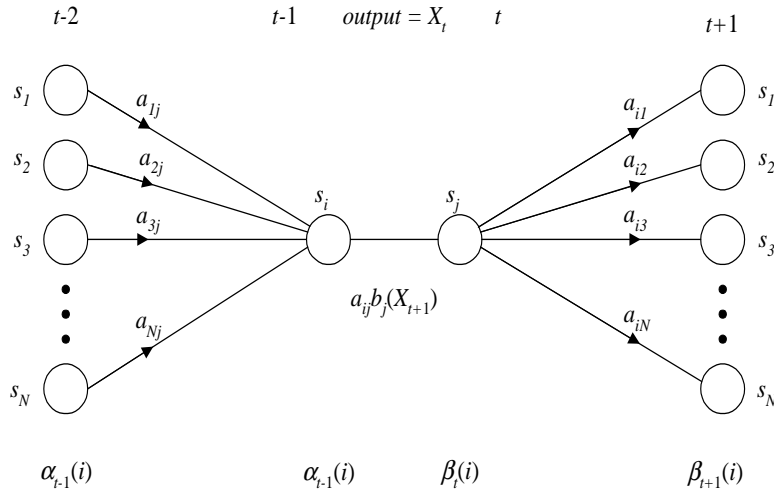


Figure 8.7 Illustration of the operations required for the computation of $\gamma_t(i, j)$, which is the probability of taking the transition from state i to state j at time t .

Since we separate the Q -function into three independent terms, the maximization procedure on $Q(\Phi, \hat{\Phi})$ can be done by maximizing the individual terms separately, subject to probability constraints.

$$\sum_{j=1}^N a_{ij} = 1 \quad \forall \text{ all } i \quad (8.36)$$

$$\sum_{k=1}^M b_j(k) = 1 \quad \forall \text{ all } j \quad (8.37)$$

Moreover, all these terms in Eqs. (8.34) and (8.35) have the following form:

$$F(x) = \sum_i y_i \log x_i \quad (8.38)$$

where $\sum_i x_i = 1$

By using the Lagrange multipliers, the function above can be proved to achieve maximum value at

$$x_i = \frac{y_i}{\sum_i y_i} \quad (8.39)$$

Using this formation, we obtain the model estimate as³:

$$\hat{a}_{ij} = \frac{\frac{1}{P(\mathbf{X}|\Phi)} \sum_{t=1}^T P(\mathbf{X}, s_{t-1} = i, s_t = j | \Phi)}{\frac{1}{P(\mathbf{X}|\Phi)} \sum_{t=1}^T P(\mathbf{X}, s_{t-1} = i | \Phi)} = \frac{\sum_{t=1}^T \gamma_t(i, j)}{\sum_{t=1}^T \sum_{k=1}^N \gamma_t(i, k)} \quad (8.40)$$

$$\hat{b}_j(k) = \frac{\frac{1}{P(\mathbf{X}|\Phi)} \sum_{t=1}^T P(\mathbf{X}, s_t = j | \Phi) \delta(X_t, o_k)}{\frac{1}{P(\mathbf{X}|\Phi)} \sum_{t=1}^T P(\mathbf{X}, s_t = j | \Phi)} = \frac{\sum_{t \in X_t = o_k} \sum_i \gamma_t(i, j)}{\sum_{t=1}^T \sum_i \gamma_t(i, j)} \quad (8.41)$$

By carefully examining the HMM re-estimation Eqs. (8.40) and (8.41), you can see that Eq. (8.40) is essentially the ratio between the expected number of transition from state i to state j and the expected number of transitions from state i . For the output probability re-estimation Eq. (8.41), the numerator is the expected number of times the observation data

³ Notice that the initial probability $\hat{\pi}_i$ can be derived as a special case of the transition probability. $\hat{\pi}_i$ is often fixed (i.e., $\hat{\pi}_i = 1$ for the initial state) for most speech applications.

emitted from state j with the observation symbol o_k , and the denominator is the expected number of times the observation data emitted from state j .

According to the EM algorithm, the forward-backward (Baum-Welch) algorithm guarantees a monotonic likelihood improvement on each iteration, and eventually the likelihood converges to a local maximum. The forward-backward algorithm can be described in a way similar to the general EM algorithm as shown in Algorithm 8.4.

ALGORITHM 8.4 THE FORWARD-BACKWARD ALGORITHM

- Step 1: Initialization: Choose an initial estimate Φ .
- Step 2: E-step: Compute auxiliary function $Q(\Phi, \hat{\Phi})$ based on Φ .
- Step 3: M-step: Compute $\hat{\Phi}$ according to the re-estimation Eqs. (8.40) and (8.41) to maximize the auxiliary Q-function.
- Step 4: Iteration: Set $\Phi = \hat{\Phi}$, repeat from step 2 until convergence.

Although the forward-backward algorithm described above is based on one training observation sequence, it can be easily generalized to multiple training observation sequences under the independence assumption between these sequences. To train an HMM from M data sequences is equivalent to finding the HMM parameter vector Φ that maximizes the joint likelihood:

$$\prod_{i=1}^M P(\mathbf{X}_i | \Phi) \quad (8.42)$$

The training procedure performs the forward-backward algorithm on each independent observation sequence to calculate the expectations (or sometimes referred to as counts) in Eqs. (8.40) and (8.41). These counts in the denominator and numerator, respectively, can be added across M data sequences respectively. Finally, all the model parameters (probabilities) are normalized to make them sum up to one. This constitutes one iteration of Baum-Welch re-estimation; iteration continues until convergence. This procedure is practical and useful because it allows you to train a good HMM in a typical speech recognition scenario where a large amount of training data is available.

For example, we let $\gamma_i^m(i, j)$ denote the $\gamma_i(i, j)$ from the m th data sequence and T^m denote the corresponding length, Equation (8.40) can be extended as:

$$\hat{a}_{ij} = \frac{\sum_{m=1}^M \sum_{t=1}^{T^m} \gamma_i^m(i, j)}{\sum_{m=1}^M \sum_{t=1}^{T^m} \sum_{k=1}^N \gamma_i^m(i, k)} \quad (8.43)$$

8.3. CONTINUOUS AND SEMI-CONTINUOUS HMMs

If the observation does not come from a finite set, but from a continuous space, the discrete output distribution discussed in the previous sections needs to be modified. The difference between the discrete and the continuous HMM lies in a different form of output probability functions. For speech recognition, use of continuous HMMs implies that the quantization procedure to map observation vectors from the continuous space to the discrete space for the discrete HMM is no longer necessary. Thus, the inherent quantization error can be eliminated.

8.3.1. Continuous Mixture Density HMMs

In choosing continuous output probability density functions $b_j(\mathbf{x})$, the first candidate is multivariate Gaussian mixture density functions. This is because they can approximate any continuous probability density function, as discussed in Chapter 3. With M Gaussian mixture density functions, we have:

$$b_j(\mathbf{x}) = \sum_{k=1}^M c_{jk} N(\mathbf{x}, \boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk}) = \sum_{k=1}^M c_{jk} b_{jk}(\mathbf{x}) \quad (8.44)$$

where $N(\mathbf{x}, \boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk})$ or $b_{jk}(\mathbf{x})$ denotes a single Gaussian density function with mean vector $\boldsymbol{\mu}_{jk}$ and covariance matrix $\boldsymbol{\Sigma}_{jk}$ for state j , M denotes the number of mixture-components, and c_{jk} is the weight for the k th mixture component satisfying

$$\sum_{k=1}^M c_{jk} = 1 \quad (8.45)$$

To take the same *divide and conquer* approach as Eq. (8.33), we need to express $b_j(\mathbf{x})$ with respect to each single mixture component as:

$$\begin{aligned} p(\mathbf{X}, \mathbf{S} | \Phi) &= \prod_{t=1}^T a_{s_{t-1}s_t} b_{s_t}(\mathbf{x}_t) \\ &= \sum_{k_1=1}^M \sum_{k_2=1}^M \dots \sum_{k_T=1}^M \left\{ \prod_{t=1}^T a_{s_{t-1}s_t} b_{s_t k_t}(\mathbf{x}_t) c_{s_t k_t} \right\} \end{aligned} \quad (8.46)$$

In the summand of Eq. (8.46) it can be considered as the summation of density with all the possible state sequences \mathbf{S} and all the possible mixture components \mathbf{K} , defined in Ω^T as the T th Cartesian product of $\Omega = \{1, 2, \dots, M\}$, as follows:

$$p(\mathbf{X}, \mathbf{S}, \mathbf{K} | \Phi) = \prod_{t=1}^T a_{s_{t-1}s_t} b_{s_t k_t}(\mathbf{x}_t) c_{s_t k_t} \quad (8.47)$$

Therefore, the joint probability density is

$$p(\mathbf{X} | \Phi) = \sum_{\mathbf{S}} \sum_{\mathbf{K} \in \Omega^T} p(\mathbf{X}, \mathbf{S}, \mathbf{K} | \Phi) \quad (8.48)$$

An auxiliary function $Q(\Phi, \hat{\Phi})$ of two model points, Φ and $\hat{\Phi}$, given an observation \mathbf{X} , can be written as:

$$Q(\Phi, \hat{\Phi}) = \sum_{\mathbf{S}} \sum_{\mathbf{K} \in \Omega^T} \frac{p(\mathbf{X}, \mathbf{S}, \mathbf{K} | \Phi)}{p(\mathbf{X} | \Phi)} \log p(\mathbf{X}, \mathbf{S}, \mathbf{K} | \hat{\Phi}) \quad (8.49)$$

From (8.47), the following decomposition can be shown:

$$\begin{aligned} & \log p(\mathbf{X}, \mathbf{S}, \mathbf{K} | \hat{\Phi}) \\ &= \sum_{t=1}^T \log \hat{a}_{s_{t-1}s_t} + \sum_{t=1}^T \log \hat{b}_{s_t k_t}(\mathbf{x}_t) + \sum_{t=1}^T \log \hat{c}_{s_t k_t} \end{aligned} \quad (8.50)$$

Maximization of the likelihood by way of re-estimation can be accomplished on individual parameter sets owing to the separability shown in (8.50). The separation of (8.50) is the key to the increased versatility of a re-estimation algorithm in accommodating mixture observation densities. The auxiliary function can be rewritten in a separated form in a similar manner as Eq. (8.33):

$$\begin{aligned} Q(\Phi, \hat{\Phi}) &= \sum_{\mathbf{S}} \sum_{\mathbf{K}} \frac{p(\mathbf{X}, \mathbf{S}, \mathbf{K} | \Phi)}{p(\mathbf{X} | \Phi)} \log p(\mathbf{X}, \mathbf{S}, \mathbf{K} | \hat{\Phi}) \\ &= \sum_{i=1}^N Q_{a_i}(\Phi, \hat{a}_i) + \sum_{j=1}^N \sum_{k=1}^M Q_{b_{jk}}(\Phi, \hat{b}_{jk}) + \sum_{j=1}^N \sum_{k=1}^M Q_{c_{jk}}(\Phi, \hat{c}_{jk}) \end{aligned} \quad (8.51)$$

The only difference we have is:

$$Q_{b_{jk}}(\Phi, \hat{b}_{jk}) = \sum_{t=1}^T p(s_t = j, k_t = k | \mathbf{X}, \Phi) \log \hat{b}_{jk}(\mathbf{x}_t), \quad (8.52)$$

and

$$Q_{c_{jk}}(\Phi, \hat{c}_{jk}) = \sum_{t=1}^T p(s_t = j, k_t = k | \mathbf{X}, \Phi) \log \hat{c}_{jk} \quad (8.53)$$

The optimization procedure is similar to what is discussed in the discrete HMM. The only major difference is $Q_{b_{jk}}(\Phi, \hat{b}_{jk})$. Maximization of $Q_{b_{jk}}(\Phi, \hat{b}_{jk})$ with respect to \hat{b}_{jk} is obtained through differentiation with respect to $\{\boldsymbol{\mu}_{jk}, \boldsymbol{\Sigma}_{jk}\}$ that satisfies:

$$\nabla_{\hat{b}_{jk}} Q_{b_{jk}}(\Phi, \hat{b}_{jk}) = 0 \quad (8.54)$$

The solutions are:

$$\hat{\boldsymbol{\mu}}_{jk} = \frac{\sum_{t=1}^T \frac{p(\mathbf{X}, s_t = j, k_t = k | \Phi)}{p(\mathbf{X} | \Phi)} \mathbf{x}_t}{\sum_{t=1}^T \frac{p(\mathbf{X}, s_t = j, k_t = k | \Phi)}{p(\mathbf{X} | \Phi)}} = \frac{\sum_{t=1}^T \zeta_t(j, k) \mathbf{x}_t}{\sum_{t=1}^T \zeta_t(j, k)} \quad (8.55)$$

$$\begin{aligned} \hat{\boldsymbol{\Sigma}}_{jk} &= \frac{\sum_{t=1}^T \frac{p(X, s_t = j, k_t = k | \Phi)}{p(\mathbf{X} | \Phi)} (\mathbf{x}_t - \hat{\boldsymbol{\mu}}_{jk})(\mathbf{x}_t - \hat{\boldsymbol{\mu}}_{jk})^t}{\sum_{t=1}^T \frac{p(X, s_t = j, k_t = k | \Phi)}{p(\mathbf{X} | \Phi)}} \\ &= \frac{\sum_{t=1}^T \zeta_t(j, k) (\mathbf{x}_t - \hat{\boldsymbol{\mu}}_{jk})(\mathbf{x}_t - \hat{\boldsymbol{\mu}}_{jk})^t}{\sum_{t=1}^T \zeta_t(j, k)} \end{aligned} \quad (8.56)$$

where $\zeta_t(j, k)$ is computed as:

$$\zeta_t(j, k) = \frac{p(\mathbf{X}, s_t = j, k_t = k | \Phi)}{p(\mathbf{X} | \Phi)} = \frac{\sum_i \alpha_{t-1}(i) a_{ij} c_{jk} b_{jk}(\mathbf{x}_t) \beta_t(j)}{\sum_{i=1}^N \alpha_t(i)} \quad (8.57)$$

In a similar manner to the discrete HMM, we can derive reestimation equation for c_{jk} as follows:

$$\hat{c}_{jk} = \frac{\sum_{t=1}^T \zeta_t(j, k)}{\sum_{t=1}^T \sum_k \zeta_t(j, k)} \quad (8.58)$$

8.3.2. Semi-continuous HMMs

Traditionally, the discrete and the continuous mixture density HMMs have been treated separately. In fact, the gap between them can be bridged under some minor assumptions with the so-called semi-continuous or tied-mixture HMM. It assumes the mixture density functions are tied together across all the models to form a set of shared kernels. In the discrete HMM, the VQ codebook is typically used to map the continuous input feature vector \mathbf{x} to o_k , so we can use the discrete output probability distribution $b_j(k)$. The codebook can be essentially regarded as one of such shared kernels. Accordingly, Eq. (8.44) can be modified as:

$$b_j(\mathbf{x}) = \sum_{k=1}^M b_j(k) f(\mathbf{x} | o_k) = \sum_{k=1}^M b_j(k) N(\mathbf{x}, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (8.59)$$

where o_k is the k th codeword, $b_j(k)$ is the same output probability distribution in the discrete HMM or the mixture weights for the continuous mixture density function, and $N(\mathbf{x}, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ are assumed to be independent of the Markov model and they are shared across all the Markov models with a very large number of mixtures M .

From the discrete HMM point of view, the needed VQ codebook consists of M continuous probability density functions, and each *codeword* has a mean vector and a covariance matrix. Typical quantization produces a codeword index that has minimum distortion to the given continuous observation \mathbf{x} . In the semi-continuous HMM, the quantization operation produces values of continuous probability density functions $f(\mathbf{x} | o_k)$ for all the codewords o_k . The structure of the semi-continuous model can be roughly the same as that of the discrete one. However, the output probabilities are no longer used directly as in the discrete HMM. In contrast, the VQ codebook density functions, $N(\mathbf{x}, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$, are combined with the discrete output probability as in Eq. (8.59). This is a combination of *discrete* model-dependent weighting coefficients with the *continuous* codebook probability density functions. Such a representation can be used to re-estimate the original VQ codebook together with the HMM.

The semi-continuous model also resembles the M -mixture continuous HMM with all the continuous output probability density functions shared among all Markov states. Compared with the continuous mixture HMM, the semi-continuous HMM can maintain the modeling ability of large-mixture probability density functions. In addition, the number of free parameters and the computational complexity can be reduced, because all the probability density functions are tied together, thus providing a good compromise between detailed acoustic modeling and trainability.

In practice, because M is large, Eq. (8.59) can be simplified by using the L most significant values $f(\mathbf{x} | o_k)$ for each \mathbf{x} without affecting the performance. Experience has shown that values of L in the range of 1-3% of M are adequate. This can be conveniently obtained during the VQ operations by sorting the VQ output and keeping the L most significant values. Let $\eta(\mathbf{x})$ denote the set of L VQ codewords that has the largest $f(\mathbf{x} | o_k)$ for the given \mathbf{x} . Then we have:

$$b_j(\mathbf{x}) \cong \sum_{k \in \eta(\mathbf{x})} p(\mathbf{x} | o_k) b_j(k) \quad (8.60)$$

Since the number of mixture components in $\eta(\mathbf{x})$ is of lower order than M , Eq. (8.60) can significantly reduce the amount of computation. In fact, $\eta(\mathbf{x})$ is the key to bridge the gap between the continuous and discrete HMM. If $\eta(\mathbf{x})$ contains only the most significant $p(\mathbf{x} | v_k)$ (i.e., only the closest codeword to \mathbf{x}), the semi-continuous HMM degenerates to the discrete HMM. On the other hand, a large VQ codebook can be used such that each Markov state could contain a number of its own codewords (a mixture of probability density

functions). The discrete output probability $b_{ij}(k)$ thus becomes the mixture weights for each state. This would go to the other extreme, a standard continuous mixture density model. We can also define $\eta(\mathbf{x})$ in such a way that we can have partial tying of $f(\mathbf{x}|o_k)$ for different phonetic classes. For example, we can have a phone-dependent codebook.

When we have a tied VQ codebook, re-estimation of these mean vectors and covariance matrices of different models will involve interdependencies. If any observation \mathbf{x}_t (no matter what model it is designated for) has a large value of posterior probability $\zeta_t(j, k)$, it will have a large contribution on re-estimation of parameters of codeword o_k . We can compute the posterior probability for each codeword from $\zeta_t(j, k)$ as defined in Eq. (8.57).

$$\zeta_t(k) = p(\mathbf{x}_t = o_k | \mathbf{X}, \Phi) = \sum_j \zeta_t(j, k) \quad (8.61)$$

The re-estimation formulas for the tied mixture can be written as:

$$\hat{\boldsymbol{\mu}}_k = \frac{\sum_{t=1}^T \zeta_t(k) \mathbf{x}_t}{\sum_{t=1}^T \zeta_t(k)} \quad (8.62)$$

$$\hat{\Sigma}_k = \frac{\sum_{t=1}^T \zeta_t(k) (\mathbf{x}_t - \hat{\boldsymbol{\mu}}_k)(\mathbf{x}_t - \hat{\boldsymbol{\mu}}_k)^t}{\sum_{t=1}^T \zeta_t(k)} \quad (8.63)$$

8.4. PRACTICAL ISSUES IN USING HMMs

While the HMM provides a solid framework for speech modeling, there are a number of issues you need to understand to make effective use of spoken language processing. In this section we point out some of the key issues related to practical applications. For expedience, we mostly use the discrete HMM as our example here.

8.4.1. Initial Estimates

In theory, the re-estimation algorithm of the HMM should reach a local maximum for the likelihood function. A key question is how to choose the right initial estimates of the HMM parameters so that the local maximum becomes the global maximum.

In the discrete HMM, if a probability is initialized to be zero, it will remain zero forever. Thus, it is important to have a reasonable set of initial estimates. Empirical study has shown that, for discrete HMMs, you can use a uniform distribution as the initial estimate. It

works reasonably well for most speech applications, though good initial estimates are always helpful to compute the output probabilities.

If continuous mixture density HMMs are used, good initial estimates for the Gaussian density functions are essential. There are a number of ways to obtain such initial estimates:

- You can use the k -means clustering procedure, as used in vector quantization clustering. The Markov state segmentation can be derived from the discrete HMM, since it is not very sensitive to the initial parameters. Based on the segmented data, you can use the k -means algorithm to derive needed Gaussian mean and covariance parameters. The mixture coefficients can be based on the uniform distribution.
- You can estimate the semi-continuous HMM from the discrete HMM. You simply need to estimate an additional covariance matrix for each VQ codeword and run an additional four or five iterations to refine the semi-continuous HMM based on the discrete HMM, which typically requires four or five iterations from the uniform distribution. When the semi-continuous HMM is trained, you take the top M *codewords*, and in each Markov state use them as the initial Gaussian density functions for the continuous density mixture model.
- You can start training a single mixture Gaussian model. You can compute the parameters from previously segmented data. You can then iteratively split the Gaussian density function in a way similar to VQ codebook generation. You typically need two or three iterations to refine the continuous density after each splitting.

8.4.2. Model Topology

Speech is a time-evolving nonstationary signal. Each HMM state has the ability to capture some quasi-stationary segment in the non-stationary speech signal. A left-to-right topology, as illustrated in Figure 8.8, is a natural candidate to model the speech signal. It has a self-transition to each state that can be used to model contiguous speech features belonging to the same state. When the quasi-stationary speech segment evolves, the left-to-right transition enables a natural progression of such evolution. In such a topology, each state has a state-dependent output probability distribution that can be used to interpret the observable speech signal. This topology is, in fact, one of the most popular HMM structures used in state-of-the-art speech recognition systems.

The state-dependent output probability distribution can be either discrete distributions or a mixture of continuous density functions. This is a special case of the transition-dependent output probability distributions. The state-dependent output probabilities can be regarded as if the transition arch-dependent output probability distributions were tied together for each state.

For the state-dependent left-to-right HMM, the most important parameter in determining the topology is the number of states. The choice of model topology depends on available

training data and what the model is used for. If each HMM is used to represent a phone, you need to have at least three to five output distributions. If such a model is used to represent a word, more states are generally required, depending on the pronunciation and duration of the word. For example, the word *tetrahydrocannabino* should have a large number of states in comparison to the word *a*. You may use at least 24 states for the former and three states for the latter. If you have the number of states depending on the duration of the signal, you may want to use 15 to 25 states for each second of speech signal. One exception is that, for silence, you may want to have a simpler topology. This is because silence is stationary, and one or two states will be sufficient.

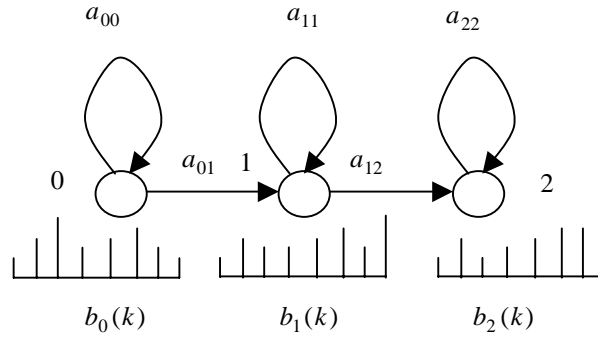


Figure 8.8 A typical hidden Markov model used to model phoneme. There are three states (0-2) and each state has an associated output probability distribution.

In practice, it is convenient to define a *null transition*. This is convenient if we want to simply traverse the HMM without consuming any observation symbol. To incorporate the null arc, you need to slightly modify the basic forward-backward or Viterbi probability equations, provided that no loops of empty transitions exist. If we denote the empty transition between state i and j as a_{ij}^e , they need to satisfy the following constraints:

$$\sum_j a_{ij} + a_{ij}^e = 1, \forall i \quad (8.64)$$

The forward probability can be augmented as follows:

$$\alpha_t(j) = \left[\sum_{i=1}^N \alpha_{t-1}(i) a_{ij} b_i(\mathbf{x}_t) + \sum_{i=1}^N \alpha_t(i) a_{ij}^e \right] \quad 1 \leq t \leq T; \quad 1 \leq j \leq N \quad (8.65)$$

Equation (8.65) appears to have a recursion, but it actually uses the value of the same time column of $\alpha_t(i)$, provided that i is already computed, which is easily achievable if we have a left-to-right empty transitions without loops of empty transitions.

8.4.3. Training Criteria

The argument for maximum likelihood estimation (MLE) is based on an assumption that the true distribution of speech is a member of the family of distributions used. This amounts to the assertion that the observed speech is genuinely produced by the HMM being used, and the only unknown parameters are the values. However, this can be challenged. Typical HMMs make many inaccurate assumptions about the speech production process, such as the output-independence assumption, the Markov assumption, and the continuous probability density assumption. Such inaccurate assumptions substantially weaken the rationale for maximum likelihood criteria. For instance, although maximum likelihood estimation is consistent (convergence to the true value), it is meaningless to have such a property if the wrong model is used. The true parameters in such cases will be the true parameters of the wrong models. Therefore, an estimation criterion that can work well in spite of these inaccurate assumptions should offer improved recognition accuracy compared with the maximum likelihood criterion. These alternative criteria include the MCE and MMIE, as discussed in Chapter 4. Finally, if we have prior knowledge about the model distribution, we can employ the Bayes method such as MAP that can effectively combine both the prior and posterior distributions in a consistent way, which is particularly suitable for adaptation or dealing with insufficient training data.

Among all these criteria, MLE remains one of the most widely used, because of its simplicity and superior performance when appropriate assumptions are made about the system design. MCE and MMIE work well for small- to medium-vocabulary speech recognition [2, 26, 36]. You can train a number of other iterations based on the ML estimates. Neither MCE nor MMIE has been found extremely effective for large-vocabulary speech recognition. However, it is possible to combine the MMIE or MCE model with the MLE model for improved performance. This is because the error patterns generated from these different models are not the same. We can decode the test utterance with these different models and vote for the most consistent results [15, 25]. We discuss MAP methods in Chapter 9, since it is mostly helpful for speaker adaptive speech recognition.

8.4.4. Deleted Interpolation

For improved robustness, it is often necessary to combine well-trained general models (such as speaker-independent) with those that are less well-trained but more detailed (such as speaker-dependent). For example, you can typically improve speech recognition accuracy with speaker-dependent training. Nevertheless, you may not have sufficient data for a particular speaker so it is desirable to use a speaker-independent model that is more general but less accurate to smooth the speaker-dependent model. One effective way to achieve robustness is to interpolate both models with a technique called deleted interpolation, in which the interpolation weights are estimated using cross-validation data. The objective function is to maximize the probability of the model generating the held-out data.

Now, let us assume that we want to interpolate two sets of models $[P_A(\mathbf{x}) \text{ and } P_B(\mathbf{x})]$, which can be either discrete probability distributions or continuous density functions] to form an interpolated model $P_{DI}(\mathbf{x})$. The interpolation procedure can be expressed as follows:

$$P_{DI}(\mathbf{x}) = \lambda P_A(\mathbf{x}) + (1 - \lambda) P_B(\mathbf{x}) \quad (8.66)$$

where the interpolation weight λ is what we need to derive from the training data.

Consider that we want to interpolate a speaker-independent model $P_A(\mathbf{x})$ with a speaker-dependent model $P_B(\mathbf{x})$. If we use speaker-independent data to estimate the interpolation weight, we may not capture needed speaker-specific information that should be reflected in the interpolation weights. What is worse is that the interpolation weight for the speaker-independent model should be equal to 1.0 if we use the same speaker-independent data from which the model was derived to estimate the interpolation weight. This is because of the MLE criterion. If we use speaker-dependent data instead, we have the weight for the speaker-dependent model equal 1.0 without achieving the desired smoothing effect. Thus the interpolation weights need to be trained using different data or *deleted* data with the so called cross-validation method.

ALGORITHM 8.5 DELETED INTERPOLATION PROCEDURE

Step 1: Initialize λ with a guessed estimate.

Step 2: Update λ by the following formula:

$$\hat{\lambda} = \frac{1}{M} \sum_{j=1}^M \sum_{t=1}^{n_j} \frac{\lambda P_{A-j}(\mathbf{x}_t^j)}{\lambda P_{A-j}(\mathbf{x}_t^j) + (1 - \lambda) P_{B-j}(\mathbf{x}_t^j)} \quad (8.67)$$

where $P_{A-j}(\mathbf{x}_t^j)$ and $P_{B-j}(\mathbf{x}_t^j)$ is $P_A(\mathbf{x})$ and $P_B(\mathbf{x})$ estimated by the entire training corpus except part j , the deleted part, respectively.; n_j is the total number of data points in part j that have been aligned to estimate the model; and \mathbf{x}_t^j indicates the t -th data point in the j -th set of the aligned data.

Step 3: If the new value $\hat{\lambda}$ is sufficiently close to the previous value λ , stop. Otherwise, go to Step 2.

We can have the training data normally divided into M parts, and train a set of $P_A(\mathbf{x})$ and $P_B(\mathbf{x})$ models using the standard EM algorithm from each combination of $M-1$ parts, with the deleted part serving as the unseen data to estimate the interpolation weights λ . These M sets of interpolation weights are then averaged to obtain the final weights.

In fact, the interpolation weights in Eq. (8.66) are similar to the Gaussian mixture weights, although $P_A(\mathbf{x})$ and $P_B(\mathbf{x})$ may not be Gaussian density functions. When we have M sets of data, we can use the same EM algorithm to estimate the interpolation weights as illustrated in Algorithm 8.5.

The deleted interpolation procedure described above can be applied after each training iteration. Then, for the following iteration of training, the learned interpolation weights can be used as illustrated in Eq. (8.66) to compute the forward-backward paths or the Viterbi maximum path. We can also have more than two distributions interpolated together. Deleted interpolation has been widely used in both acoustic and language modeling where smoothing is needed.

8.4.5. Parameter Smoothing

One simple reality for probabilistic modeling is that as many observations as possible are required to reliably estimate model parameters. However, in reality, only a finite amount of training data is available. If the training data are limited, this will result in some parameters being inadequately trained, and classification based on poorly trained models will result in higher recognition error rate. There are many possible solutions to address the problem of insufficient training data:

- You can increase the size of the training data. *There is no data like more data.*
- You can reduce the number of free parameters to be re-estimated. This has its limitations, because a number of significant parameters are always needed to model physical events.
- You can interpolate one set of parameter estimates with another set of parameter estimates, for which an adequate amount of training data exists. Deleted interpolation discussed in Section 8.4.4, can be used effectively. In the discrete HMM, one simple approach is to set a floor to both the transition probability and the output probability in order to eliminate possible zero estimates. The same principle applies to the SCHMM as well as the mixing coefficients of the continuous density HMM. Parameter flooring can be regarded as a special case of interpolation with the uniform distribution.
- You can tie parameters together to reduce the number of free parameters. The SCHMM is a typical example of such parameter-tying techniques.

For the continuous mixture HMM, you need to pay extra attention to smoothing the covariance matrices. There are a number of techniques you can use:

- You can interpolate the covariance matrix with those that are better trained or a priori via the MAP method.
- You can tie the Gaussian covariance matrices across different mixture components or across different Markov states. A very general shared Gaussian density model is discussed in [20].

- You can use the diagonal covariance matrices if the correlation among feature coefficients is weak, which is the case if you use uncorrelated features such as the MFCC.
- You can combine these methods together.

In practice, we can reduce the speech recognition error rate by 5-20% with various smoothing techniques, depending on the available amount of training data.

8.4.6. Probability Representations

When we compute the forward and backward probabilities in the forward-backward algorithm, they will approach zero in exponential fashion if the observation sequence length, T , becomes large enough. For sufficiently large T , the dynamic range of these probabilities will exceed the precision range of essentially any machine. Thus, in practice, it will result in underflow on the computer if probabilities are represented directly. We can resolve this implementation problem by scaling these probabilities with some scaling coefficient so that they remain within the dynamic range of the computer. All of these scaling coefficients can be removed at the end of the computation without affecting the overall precision.

For example, let $\alpha_t(i)$ be multiplied by a scaling coefficient, S_t :

$$S_t = 1 / \sum_i \alpha_t(i) \quad (8.68)$$

so that $\sum_i S_t \alpha_t(i) = 1$ for $1 \leq t \leq T$. $\beta_t(i)$ can also be multiplied by S_t for $1 \leq t \leq T$. The recursion involved in computing the forward and backward variables can be scaled at each stage of time t by S_t . Notice that $\alpha_t(i)$ and $\beta_t(i)$ are computed recursively in exponential fashion; therefore, at time t , the total scaled factor applied to the forward variable $\alpha_t(i)$ is

$$Scale_\alpha(t) = \prod_{k=1}^t S_k \quad (8.69)$$

and the total scaled factor applied to the backward variable $\beta_t(i)$ is

$$Scale_\beta(t) = \prod_{k=t}^T S_k \quad (8.70)$$

This is because the individual scaling factors are multiplied together in the forward and backward recursion. Let $\alpha'_t(i)$, $\beta'_t(i)$, and $\gamma'_t(i, j)$ denote their corresponding scaled variables, respectively. Note that

$$\sum_i \alpha'_T(i) = Scale_\alpha(T) \sum_i \alpha_T(i) = Scale_\alpha(T) P(\mathbf{X} | \Phi) \quad (8.71)$$

The scaled intermediate probability, $\gamma'_t(i, j)$, can then be written as:

$$\gamma'_t(i, j) = \frac{Scale_\alpha(t-1)\alpha_{t-1}(i)a_{ij}b_j(X_t)\beta_t(j)Scale_\beta(t)}{Scale_\alpha(T)\sum_{i=1}^N\alpha_T(i)} = \gamma_t(i, j) \quad (8.72)$$

Thus, the intermediate probabilities can be used in the same way as unscaled probabilities, because the scaling factor is cancelled out in Eq. (8.72). Therefore, re-estimation formulas can be kept exactly except that $P(\mathbf{X} | \Phi)$ should be computed according to

$$P(\mathbf{X} | \Phi) = \sum_i \alpha'_T(i) / Scale_\alpha(T) \quad (8.73)$$

In practice, the scaling operation need not be performed at every observation time. It can be used at any scaling interval for which the underflow is likely to occur. In the unscaled interval, $Scale_\alpha$ can be kept as unity.

An alternative way to avoid underflow is to use a logarithmic representation for all the probabilities. This not only ensures that scaling is unnecessary, as underflow cannot happen, but also offers the benefit that integers can be used to represent the logarithmic values, thereby changing floating point operations to fixed point ones, which is particularly suitable for Viterbi-style computation, as Eq. (8.25) requires no probability addition.

In the forward-backward algorithm we need to have probability addition. We can keep a table on $\log_b P_2 - \log_b P_1$. If we represent probability P by $\log_b P$, more precision can be obtained by setting b closer to unity. Let us assume that we want to add P_1 and P_2 and that $P_1 \geq P_2$. We have:

$$\log_b (P_1 + P_2) = \log_b P_1 + \log_b (1 + b^{\log_b P_2 - \log_b P_1}) \quad (8.74)$$

If P_2 is so many orders of magnitude smaller than P_1 , adding the two numbers will just result in P_1 . We could store all possible values of $(\log_b P_2 - \log_b P_1)$. Using logarithms introduces errors for addition operation. In practice, double-precision float representation can be used to minimize the impact of the precision problems.

8.5. HMM LIMITATIONS

There are a number of limitations in the conventional HMMs. For example, HMMs assume the duration follows an exponential distribution, the transition probability depends only on the origin and destination, and all observation frames are dependent only on the state that generated them, not on neighboring observation frames. Researchers have proposed a number of techniques to address these limitations, albeit these solutions have not significantly improved speech recognition accuracy for practical applications.

8.5.1. Duration Modeling

One major weakness of conventional HMMs is that they do not provide an adequate representation of the temporal structure of speech. This is because the probability of state occupancy decreases exponentially with time as shown in Eq. (8.75). The probability of t consecutive observations in state i is the probability of taking the self-loop at state i for t times, which can be written as

$$d_i(t) = a_{ii}^t (1 - a_{ii}) \quad (8.75)$$

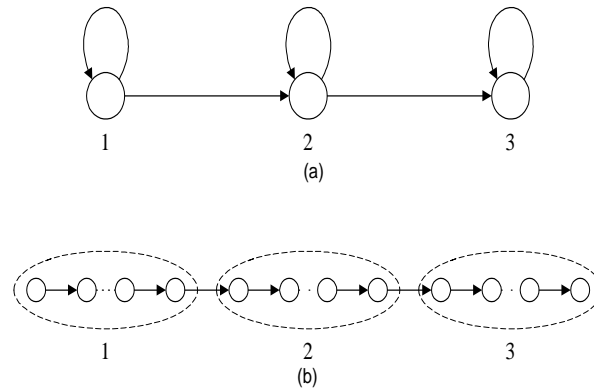


Figure 8.9 A standard HMM (a) and its corresponding explicit duration HMM (b) where the self transitions are replaced with the explicit duration probability distribution for each state.

An improvement to the standard HMM results from the use of HMMs with an explicit time duration distribution for each state [30, 39]. To explain the principle of time duration modeling, a conventional HMM with exponential state duration density and a time duration HMM with specified state duration densities (which can be either a discrete distribution or a continuous density) are illustrated in Figure 8.9. In (a), the state duration probability has an exponential form as in Eq. (8.75). In (b), the self-transition probabilities are replaced with an explicit duration probability distribution. At time t , the process enters state i for duration τ with probability density $d_i(\tau)$, during which the observations $X_{t+1}, X_{t+2}, \dots, X_{t+\tau}$ are generated. It then transfers to state j with transition probability a_{ij} only after the appropriate τ observations have occurred in state i . Thus, by setting the time duration probability density to be the exponential density of Eq. (8.75) the time duration HMM can be made equivalent to the standard HMM. The parameters $d_i(\tau)$ can be estimated from observations along with the other parameters of the HMM. For expedience, the duration density is usually truncated at a maximum duration value T_d . To re-estimate the parameters of the HMM with time duration modeling, the forward recursion must be modified as follows:

$$\alpha_t(j) = \sum_{\tau} \sum_{i, i \neq j} \alpha_{t-\tau}(i) a_{ij} d_j(\tau) \prod_{l=1}^{\tau} b_j(X_{t-\tau+l}) \quad (8.76)$$

where the transition from state i to state j depends not only upon the transition probability a_{ij} but also upon all the possible time periods τ that may occur in state j . Intuitively, Eq. (8.76) illustrates that when state j is reached from previous states i , the observations may stay in state j for a period of τ with duration density $d_j(\tau)$, and each observation emits its own output probability. All possible durations must be considered, which leads to summation with respect to τ . The independence assumption of observations results in the Π term of the output probabilities. Similarly, the backward recursion can be written as:

$$\beta_t(i) = \sum_{\tau} \sum_{j, j \neq i} a_{ij} d_j(\tau) \prod_{l=1}^{\tau} b_j(X_{t+l}) \beta_{t+\tau}(j) \quad (8.77)$$

The modified Baum-Welch algorithm can then be used based on Eq. (8.76) and (8.77). The proof of the re-estimation algorithm can be based on the modified Q -function except that $P(\mathbf{X}, \mathbf{S} | \Phi)$ should be replaced with $P(\mathbf{X}, \mathbf{S}, \mathbf{T} | \Phi)$, which denotes the joint probability of observation, \mathbf{X} , state sequence, $\mathbf{S} = \{s_1, s_2, \dots, s_k, \dots, s_{N_s}\}$ in terms of state s_k with time duration τ_k , and the corresponding duration sequence, $\mathbf{T} = \{\tau_1, \tau_2, \dots, \tau_k, \dots, \tau_{N_s}\}$.

$$Q(\Phi, \hat{\Phi}) = \frac{1}{P(\mathbf{X} | \Phi)} \sum_{\mathbf{T}} \sum_{\mathbf{S}} P(\mathbf{X}, \mathbf{S}, \mathbf{T} | \Phi) \log P(\mathbf{X}, \mathbf{S}, \mathbf{T} | \hat{\Phi}) \quad (8.78)$$

In a manner similar to the standard HMM, $\gamma_{t,\tau}(i, j)$ can be defined as the transition probability from state i at time t to state j with time duration τ in state j . $\gamma_{t,\tau}(i, j)$ can be written as:

$$\gamma_{t,\tau}(i, j) = \alpha_t(i) a_{ij} d_j(\tau) \prod_{l=1}^{\tau} b_j(X_{t+l}) \beta_{t+\tau}(j) / \sum_{k=1}^N \alpha_t(k) \quad (8.79)$$

Similarly, the probability of being in state j at time t with duration τ can be computed as:

$$\gamma_{t,\tau}(j) = \sum_i \gamma_{t,\tau}(i, j) \quad (8.80)$$

The re-estimation algorithm can be derived from Eq. (8.80), the Viterbi decoding algorithm can be used for the time duration model, and the optimal path can be determined according to:

$$_t(j) = \text{Max}_i \text{Max}_{\tau} [V_{t-\tau}(i) a_{ij} d_j(\tau) \prod_{l=1}^{\tau} b_j(X_{t-\tau+l})] \quad (8.81)$$

There are drawbacks to the use of the time duration modeling discussed here. One is the great increase in computational complexity by a factor of $O(D^2)$, where D is the time duration distribution length. Another problem is the large number of additional parameters D that must be estimated. One proposed remedy is to use a continuous density function instead of the discrete distribution $d_j(\tau)$.

In practice, duration models offered only modest improvement for speaker-independent continuous speech recognition. Many systems have even eliminated the transition probability completely because output probabilities are so dominant. Nevertheless, duration information is very effective for pruning unlikely candidates during the large-vocabulary speech recognition decoding process.

8.5.2. First-Order Assumption

As you can see from the previous section, the duration of each stationary segment captured by a single state is inadequately modeled. Another way to alleviate the duration problem is to eliminate the first-order transition assumption and to make the underlying state sequence a second-order Markov chain [32]. As a result, the transition probability between two states at time t depends on the states in which the process was at time $t-1$ and $t-2$. For a given state sequence $\mathbf{S} = \{s_1, s_2, \dots, s_T\}$, the probability of the state should be computed as:

$$P(\mathbf{S}) = \prod_t a_{s_{t-2}s_{t-1}s_t} \quad (8.82)$$

where $a_{s_{t-2}s_{t-1}s_t} = P(s_t | s_{t-2}s_{t-1})$ is the transition probability at time t , given the two-order state history. The re-estimation procedure can be readily extended based on Eq. (8.82). For example, the new forward probability can be re-defined as:

$$\alpha_t(j, k) = P(X_1^t, s_{t-1} = j, s_t = k | \lambda) = \sum_i \alpha_{t-1}(i, j) a_{ijk} b_k(X_t) \quad (8.83)$$

where $a_{ijk} = P(s_t = k | s_{t-2} = i, s_{t-1} = j)$. Similarly, we can define the backward probability as:

$$\beta_t(i, j) = P(X_{t+1}^T | s_{t-1} = i, s_t = j, \lambda) = \sum_k a_{ijk} b_k(X_{t+1}) \beta_{t+1}(j, k) \quad (8.84)$$

With Eq. (8.83) and (8.84), the MLE estimates can be derived easily based on the modified $\gamma_t(i, j, k)$:

$$\begin{aligned} \gamma_t(i, j, k) &= P(s_{t-1} = i, s_t = j, s_{t+1} = k, \mathbf{X} | \Phi) \\ &= \alpha_t(i, j) a_{ijk} b_k(X_{t+1}) \beta_{t+1}(j, k) / P(\mathbf{X} | \Phi) \end{aligned} \quad (8.85)$$

In practice, the second-order model is computationally very expensive as we have to consider the increased state space, which can often be realized with an equivalent first-order

hidden Markov model on the two-fold product state space. It has not offered significantly improved accuracy to justify its increase in computational complexity for most applications.

8.5.3. Conditional Independence Assumption

The third major weakness in HMMs is that all observation frames are dependent only on the state that generated them, not on neighboring observation frames. The conditional independence assumption makes it hard to effectively handle nonstationary frames that are strongly correlated. There are a number of ways to alleviate the conditional independence assumption [34]. For example, we can assume the output probability distribution depends not only on the state but also on the previous frame. Thus, the probability of a given state sequence can be rewritten as:

$$P(\mathbf{X}|\mathbf{S}, \Phi) = \prod_{t=1}^T P(X_t | X_{t-1}, s_t, \Phi) \quad (8.86)$$

As the parameter space becomes huge, we often need to quantize X_{t-1} into a smaller set of codewords so that we can keep the number of free parameters under control. Thus, Eq. (8.86) can be simplified as:

$$P(\mathbf{X} | \mathbf{S}, \Phi) = \prod_{t=1}^T P(X_t | \mathfrak{R}(X_{t-1}), s_t, \Phi) \quad (8.87)$$

where $\mathfrak{R}()$ denotes the quantized vector that has a small codebook size, L . Although this can dramatically reduce the space of the free conditional output probability distributions, the total number of free parameters will still increase by L times.

The re-estimation for conditional dependent HMMs can be derived with the modified Q -function, as discussed in the previous sections. In practice, it has not demonstrated convincing accuracy improvement for large-vocabulary speech recognition.

8.6. HISTORICAL PERSPECTIVE AND FURTHER READING

The Markov chain was named after Russian scientist A. Markov for his pioneering work in analyzing the letter sequence in the text of a literary work in 1913 [33]. In the 1960s, Baum and others further developed efficient methods for training the model parameters [4, 5]. When the observation is real valued, the use of continuous or semi-continuous HMMs can improve the overall performance. Baum et al. also developed the method to use continuous density functions that are strictly log concave [5], which was relaxed by Liporace [31] and expanded by Juang to include mixture density functions [27].

The Viterbi algorithm shares the same concept that was independently discovered by researchers in many separate fields [28], including Vintsyuk [42], Needleman and Wunsch [35], Sankoff [41], Sakoe and Chiba [40], and Wagner and Fischer [44].

Jim Baker did his Ph.D. thesis under Raj Reddy at Carnegie Mellon using HMMs for speech recognition [3]. At the same time Fred Jelinek and his colleagues at IBM Research pioneered widespread applications [23]. Since the 1980s, partly because of the DARPA-funded speech projects, HMMs have become a mainstream technique for modeling speech, as exemplified by advanced systems developed at BBN, Bell Labs, Carnegie Mellon, IBM, Microsoft, SRI, and others [9, 17, 29, 46]. The Ph.D. theses from Kai-Fu Lee [29], Hsiao-Wuen Hon [16], and Mei-Yuh Hwang [22] at Carnegie Mellon addressed many important practical issues in using HMMs for speech recognition. There are also a number of good books on the practical use of HMMs [18, 24, 38, 45].

The choice of different output probabilities depends on a number of factors such as the availability of training data, the feature characteristics, the computational complexity, and the number of free parameters [19] [34]. The semi-continuous model, also known as the tied-mixture model, was independently proposed by Huang and Jack [21] and Bellegarda and Nahamoo [6]. Other improvements include explicit duration modeling [1, 11, 13, 14, 30, 39], high-order and conditional models [7, 32, 34], which have yet to be shown effective for practical speech recognition.

Both Carnegie Mellon University's open speech software⁴ and Cambridge University's HTK⁵ are a good starting point for those interested in using the existing tools for running experiments.

The HMM have become the most prominent techniques for speech recognition today. Most of the state-of-the-art speech recognition systems on the market are based on HMMs described in this chapter.

⁴ <http://www.speech.cs.cmu.edu/sphinx/>

⁵ <http://htk.eng.cam.ac.uk/>

REFERENCES

- [1] Anastasakos, A., R. Schwartz, and H. Sun, "Duration Modeling in Large Vocabulary Speech Recognition" in *Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing* 1995, Detroit, MI, pp. 628-631.
- [2] Bahl, L.R., *et al.*, "Speech Recognition with Continuous-Parameter Hidden Markov Models," *Computer Speech and Language*, 1987, **2**, pp. 219-234.
- [3] Baker, J.K., "The DRAGON System - An Overview," *Trans. on Acoustics, Speech and Signal Processing*, 1975, **23**(1), pp. 24-29.
- [4] Baum, L.E. and J.A. Eagon, "An Inequality with Applications to Statistical Estimation for Probabilistic Functions of Markov Processes and to a Model for Ecology," *Bulletin of American Mathematical Society*, 1967, **73**, pp. 360-363.
- [5] Baum, L.E., *et al.*, "A Maximization Technique Occurring in the Statistical Analysis of Probabilistic Functions of Markov Chains," *Annals of Mathematical Statistics*, 1970, **41**, pp. 164-171.
- [6] Bellegarda, J.R. and D. Nahamoo, "Tied Mixture Continuous Parameter Models for Large Vocabulary Isolated Speech Recognition," *Int. Conf. on Acoustics, Speech and Signal Processing*, 1989 pp. 13-16.
- [7] Brown, P., *The Acoustic-Modeling Problem in Automatic Speech Recognition*, Ph.D. Thesis in *Computer Science* 1987, Carnegie Mellon University, Pittsburgh.
- [8] Brown, P.F., *et al.*, "The Mathematics of Statistical Machine Translation: Parameter Estimation," *Computational Linguistics*, 1995, **19**(2), pp. 263--312.
- [9] Chou, W., C.H. Lee, and B.H. Juang, "Minimum Error Rate Training of Inter-Word Context Dependent Acoustic Model Units in Speech Recognition" in *Proc. of the Int. Conf. on Spoken Language Processing* 1994, Yokohama, Japan, pp. 439-442.
- [10] Church, K., "A Stochastic Parts Program and Noun Phrase Parser for Unrestricted Text," *Proc. of the Second Conf. on Applied Natural Language Processing*, 1988, Austin, Texas pp. 136-143.
- [11] Deng, L., M. Lennig, and P. Mermelstein, "Use of Vowel Duration Information in a Large Vocabulary Word Recognizer," *Journal of the Acoustical Society of America*, 1989, **86**(2 August), pp. 540-548.
- [12] DeRose, S.J., "Grammatical Category Disambiguation by Statistical Optimization," *Computational Linguistics*, 1988(1), pp. 31-39.
- [13] Dumouchel, P. and D. Shaughnessy, "Segmental Duration and HMM Modeling," *Proc. of the European Conf. on Speech Communication and Technology*, 1995, Madrid, Spain pp. 803-806.
- [14] Ferguson, J.D., "Variable Duration Models for Speech" in *Proc. of the Symposium on the Application of Hidden Markov Models to Text and Speech*, J.D. Ferguson, Editor 1980, New Jersey, pp. 143-179, Princeton.
- [15] Fiscus, J., "A Post-Processing System to Yield Reduced Word Error Rates: Recognizer Output Voting Error Reduction (ROVER)," *IEEE Workshop on Automatic Speech Recognition and Understanding*, 1997, Santa Barbara, CA pp. 347-352.

- [16] Hon, H.W., *Vocabulary-Independent Speech Recognition: The VOCIND System*, Ph.D Thesis in *Department of Computer Science* 1992, Carnegie Mellon University, Pittsburgh.
- [17] Huang, X.D., *et al.*, "The SPHINX-II Speech Recognition System: An Overview," *Computer Speech and Language*, 1993 pp. 137-148.
- [18] Huang, X.D., Y. Ariki, and M.A. Jack, *Hidden Markov Models for Speech Recognition*, 1990, Edinburgh, U.K., Edinburgh University Press.
- [19] Huang, X.D., *et al.*, "A Comparative Study of Discrete, Semicontinuous, and Continuous Hidden Markov Models," *Computer Speech and Language*, 1993, 7(4), pp. 359-368.
- [20] Huang, X.D., *et al.*, "Deleted Interpolation and Density Sharing for Continuous Hidden Markov Models," *IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, 1996.
- [21] Huang, X.D. and M.A. Jack, "Semi-Continuous Hidden Markov Models with Maximum Likelihood Vector Quantization" in *IEEE Workshop on Speech Recognition* 1988.
- [22] Hwang, M., *Subphonetic Modeling in HMM-based Speech Recognition Systems*, Ph.D. Thesis in *Computer Science* 1994, Carnegie Mellon University, Pittsburgh.
- [23] Jelinek, F., "Continuous Speech Recognition by Statistical Methods," *Proc. of the IEEE*, 1976, 64(4), pp. 532-556.
- [24] Jelinek, F., *Statistical Methods for Speech Recognition*, 1998, Cambridge, MA, MIT Press.
- [25] Jiang, L. and X. Huang, "Unified Decoding and Feature Representation for Improved Speech Recognition," *Proc. of the 6th European Conf. on Speech Communication and Technology*, 1999, Budapest, Hungary pp. 1331-1334.
- [26] Juang, B.H., W. Chou, and C.H. Lee, "Statistical and Discriminative Methods for Speech Recognition" in *Automatic Speech and Speaker Recognition - Advanced Topics*, C.H. Lee, F.K. Soong, and K.K. Paliwal, eds. 1996, Boston, pp. 109-132, Kluwer Academic Publishers.
- [27] Juang, B.H., S.E. Levinson, and M.M. Sondhi, "Maximum Likelihood Estimation for Multivariate Mixture Observations of Markov Chains," *IEEE Trans. on Information Theory*, 1986, IT-32(2), pp. 307-309.
- [28] Kruskal, J., "An Overview of Sequence Comparison" in *Time Warps, String Edits, and Macromolecules: The Theory and Practice of Sequence Comparison*, D. Sankoff and J. Kruskal, eds. 1983, Reading, MA., pp. 1-44, Addison-Wesley.
- [29] Lee, K.F., *Large-Vocabulary Speaker-Independent Continuous Speech Recognition: The SPHINX System*, Ph.D. Thesis in *Computer Science Dept.* 1988, Carnegie Mellon University, Pittsburgh.
- [30] Levinson, S.E., "Continuously Variable Duration Hidden Markov Models for Automatic Speech Recognition," *Computer Speech and Language*, 1986, pp. 29-45.

- [31] Liporace, L.R., "Maximum Likelihood Estimation for Multivariate Observations of Markov Sources," *IEEE Trans. on Information Theory*, 1982, **28**, pp. 729-734.
- [32] Mari, J., J. Haton, and A. Kriouile, "Automatic Word Recognition Based on Second-Order Hidden Markov Models," *IEEE Trans. on Speech and Audio Processing*, 1977, **5**(1), pp. 22-25.
- [33] Markov, A.A., "An Example of Statistical Investigation in the Text of 'Eugene Onegin', Illustrating Coupling of Tests in Chains," *Proc. of the Academy of Sciences of St. Petersburg*, 1913, Russia pp. 153-162.
- [34] Ming, J. and F. Smith, "Modelling of the Interframe Dependence in an HMM Using Conditional Gaussian Mixtures," *Computer Speech and Language*, 1996, **10**(4), pp. 229-247.
- [35] Needleman, S. and C. Wunsch, "A General Method Applicable to the Search for Similarities in the Amino-acid Sequence of Two Proteins," *Journal of Molecular Biology*, 1970, **48**, pp. 443-453.
- [36] Normandin, Y., "Maximum Mutual Information Estimation of Hidden Markov Models" in *Automatic Speech and Speaker Recognition*, C.H. Lee, F.K. Soong, and K.K. Paliwal, eds. 1996, Norwell, MA, Kluwer Academic Publishers.
- [37] Rabiner, L.R., "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition," *Proc. of IEEE*, 1989, **77**(2), pp. 257-286.
- [38] Rabiner, L.R. and B.H. Juang, *Fundamentals of Speech Recognition*, Prentice Hall Signal Processing Series, eds. A.V. Oppenheim, 1993, Englewood Cliffs, NJ, Prentice-Hall.
- [39] Russell, M.J. and R.K. Moore, "Explicit Modeling of State Occupancy in Hidden Markov Models for Automatic Speech Recognition," *Int. Conf. on Acoustics, Speech and Signal Processing*, 1985 pp. 5-8.
- [40] Sakoe, H. and S. Chiba, "Dynamic Programming Algorithm Optimization for Spoken Word Recognition," *IEEE Trans. on Acoustics, Speech and Signal Processing*, 1978, **26**(1), pp. 43-49.
- [41] Sankoff, D., "Matching Sequences under Deletion-Insertion Constraints," *Proc. of the National Academy of Sciences*, 1972, **69**, pp. 4-6.
- [42] Vintsyuk, T.K., "Speech Discrimination by Dynamic Programming," *Cybernetics*, 1968, **4**(1), pp. 52-57.
- [43] Viterbi, A.J., "Error Bounds for Convolutional Codes and an Asymptotically Optimum Decoding Algorithm," *IEEE Trans. on Information Theory*, 1967, **13**(2), pp. 260-269.
- [44] Wagner, R. and M. Fischer, "The String-to-String Correction Problem," *Journal of the ACM*, 1974, **21**, pp. 168-173.
- [45] Waibel, A.H. and K.F. Lee, *Readings in Speech Recognition*, 1990, San Mateo, CA, Morgan Kaufman Publishers.
- [46] Young, S.J. and P.C. Woodland, "The Use of State Tying in Continuous Speech Recognition," *Proc. of Eurospeech*, 1993, Berlin pp. 2203-2206.